

APPLICATION OF EM ALGORITHM ON MISSING CATEGORICAL DATA
ANALYSIS

NORAINI BINTI HASAN

A report submitted in partial fulfilment of the
requirements for the award of the degree of
Master of Science (Mathematic)

Faculty of Science

Universiti Teknologi Malaysia

DECEMBER 2009

To my beloved husband, son and all my family members

ACKNOWLEDGEMENT

In preparing this thesis, I was in contact with many people, researchers, academicians, and practitioners. They have contributed towards my understanding and thoughts. In particular, I wish to express my sincere appreciation to my thesis supervisor, Assoc. Prof. Dr. Ismail b. Mohamad, for encouragement, guidance, critics and friendship. Without their continued support and interest, this thesis would have never been the same as presented here. Librarians at UTM also deserve my special thanks for their assistance in supplying the relevant literatures.

My colleagues should also be recognised for their support and the assistance provided at various occasions. Their views and tips are useful indeed. My sincere appreciation also extends to my beloved husband and son, my family and also not forgotten my in-laws, for their understanding and sacrificial. Unfortunately, it is not possible to list all of them in this limited space.

ABSTRAK

Algoritma EM merupakan salah satu daripada kaedah untuk menyelesaikan masalah berkaitan dengan data tidak lengkap berdasarkan kepada satu rangka lengkap. Algoritma EM merupakan satu pendekatan parametrik untuk mencari taksiran ML data tidak lengkap. Algoritma ini terbahagi kepada dua langkah, dimana langkah pertamanya, langkah Ekspektasi, atau lebih dikenali sebagai langkah E, mencari ekspektasi kepada loglikelihood, bersyarat kepada data yang dapat diperolehi dan anggaran terkini, $\theta^{(t)}$. Langkah kedua, langkah Pemaksimuman atau langkah M dimana ia akan memaksimumkan nilai loglikelihood untuk mencari satu anggaran parameter yang baru. Prosedur ini berlaku berselang-seli antara kedua-dua langkah ini sehingga anggaran parameter tersebut malar.

ABSTRACT

Expectation- Maximization algorithm, or in short, EM algorithm is one of the methodologies for solving incomplete data problems sequentially based on a complete framework. The EM algorithm is a parametric approach to find the Maximum Likelihood, ML parameter estimates for incomplete data. The algorithm consists of two steps. The first step is the Expectation step, better known as E-step, finds the expectation of the loglikelihood, conditional on the observed data and the current parameter estimates; say $Q(\theta)$. The second step is the Maximization step, or M-step, which maximize the loglikelihood to find new estimates of the parameters. The procedure alternates between the two steps until the parameter estimates converge to some fixed values.

TABLE OF CONTENTS

CHAPTER	TITLE	PAGE
	DECLARATION	ii
	DEDICATION	iii
	ACKNOWLEDGEMENTS	iv
	ABSTRACT	v
	ABSTRACT	vi
	TABLE OF CONTENTS	vii
	LIST OF TABLES	x
	LIST OF SYMBOLS	xiv
1	INTRODUCTION	1
1.1	Problem Statement	1
1.2	Objective Of The Study	2
1.3	Scope Of The Study	3

1.4	Significance Of The Study	3
2	LITERATURE REVIEW	4
2.1	Missing Data	5
2.1.1	Classes of Missing Data	6
2.1.1.1	Censored Data	6
2.1.1.2	Latent Variable	7
2.1.1.3	Non-Response Item	8
2.2	The Expectation-Maximization Algorithm	10
3	RESEARCH METHODOLOGY	12
3.1	Missing Data Patterns	12
3.2	General Definition of Missingness Mechanism	15
3.3	EM Theory in General	17
3.4	Incomplete Contingency Table	27
3.4.1	ML Estimation in Incomplete Contingency Table	27
3.4.2	The EM Algorithm	28
3.4.2.1	Multinomial Sampling	28
3.4.2.2	Product Multinomial Sampling	30
3.4.2.3	EM Algorithm to Determine the ML Estimates of Cell Probabilities in An Incomplete \times Contingency Table Data Missing on Both Categories	31

3.5	Chi- Squared Test	35
3.5.1	Goodness-of-fit Test	35
3.5.2	Independence Test	41
4	RESULT AND DISCUSSION	46
4.1	Data Construction	52
4.1.1	Missing Completely At Random (MCAR)	53
4.1.2	Missing At Random (MAR)	59
4.1.3	Not Missing At Random (NMAR)	64
4.1.4	The Chi- Squared Test	68
5	CONCLUSION AND RECOMMENDATION	72
5.1	Conclusion	72
5.2	Recommendation	73
	REFERENCES	75

LIST OF TABLES

TABLE NO.	TITLE	PAGE
3.1	Classification of sample units in an incomplete \times contingency table	32
3.2	Frequency distribution	36
3.3	The calculation of statistic	39
3.4	The observed frequency of category i	40
3.5	A two-way contingency table	41
3.6	A two-way dimensional contingency table of joint events	42
4.1	An example dataset of full data	47
4.1(a)	Continuous data	47
4.1(b)	Categorical data	47
4.2	An example dataset for MCAR	49
4.2(a)	Continuous data	49
4.2(b)	Categorical data	49
4.3	An example dataset for MAR	50
4.3(a)	Continuous data	50

4.3(b)	Categorical data	51
4.4	An example dataset for NMAR	51
4.4(a)	Continuous data	51
4.4(b)	Categorical data	52
4.5	Full data	53
4.6	Artificial Incomplete Data for MCAR.	54
4.6(a)	MCAR with 10% data missing	54
4.6(b)	MCAR with 20% data missing	54
4.6(c)	MCAR with 30% data missing	55
4.7	Marginal total of probabilities for MCAR with 10% data missing	55
4.8	Iteration of EM algorithm for MCAR with 10% data missing problem	57
4.9	Complete data obtained by EM algorithm for 10% MCAR problem	58
4.10	MCAR with 20% of the data are missing	58
4.10(a)	Iteration of EM algorithm.	58
4.10(b)	Complete data obtained by EM algorithm	58
4.11	MCAR with 30% of the data are missing	59
4.11(a)	Iteration of EM algorithm.	59
4.11(b)	Complete data obtained by EM algorithm	59
4.12	Artificial Incomplete Data for MAR.	60
4.12(a)	MAR with 10% data missing	60
4.12(b)	MAR with 20% data missing	61

4.12(c)	MAR with 30% data missing	61
4.13	MAR with 10% of the data are missing	61
4.13(a)	Iteration of EM algorithm.	61
4.13(b)	Complete data obtained by EM algorithm	62
4.14	MAR with 20% of the data are missing	62
4.14(a)	Iteration of EM algorithm.	62
4.14(b)	Complete data obtained by EM algorithm	62
4.15	MAR with 30% of the data are missing	63
4.15(a)	Iteration of EM algorithm.	63
4.15(b)	Complete data obtained by EM algorithm	63
4.16	Artificial Incomplete Data for NMAR.	65
4.16(a)	NMAR with 10% data missing	65
4.16(b)	NMAR with 20% data missing	65
4.16(c)	NMAR with 30% data missing	65
4.17	NMAR with 10% of the data are missing	66
4.17(a)	Iteration of EM algorithm.	66
4.17(b)	Complete data obtained by EM algorithm	66
4.18	MAR with 20% of the data are missing	66
4.18(a)	Iteration of EM algorithm.	66
4.18(b)	Complete data obtained by EM algorithm	67
4.19	MAR with 30% of the data are missing	67
4.19(a)	Iteration of EM algorithm.	67
4.19(b)	Complete data obtained by EM algorithm	67

4.20	The calculation for full data	69
4.21	The values for all cases	70

LIST OF SYMBOLS

The observed value

The missing value

Number of observations or Total counts

Estimates of

Current estimates of

The counts in cell (,)

The observed value for

The probability that an observation falls in cell (,)

() r th estimates of

Observed frequencies

Expected frequencies

CHAPTER 1

INTRODUCTION

1.1 PROBLEM STATEMENT

Incomplete table is referred to the table in which the entries or information on one or more of the categorical variables are missing, a prior zero or undetermined (Fienberg, 1980). Missing data treatment is an important data quality issue in data mining, data warehousing, and database management. Real-world data often has missing values.

The presence of missing values can cause serious problems when the data is used for reporting, information sharing, and decision support. First, data with missing values may provide biased information. For example, a survey question that is related to personal information will more likely be left unanswered for those who are more sensitive about privacy. Second, many data modeling and analysis techniques cannot deal with missing values and have to cast out a whole record value if one of the attribute values is missing. Third, even though some data modeling and analysis tools can handle missing values, there are often restrictions in the domain of missing values. For example, classification systems typically do not allow missing values in the class attribute.

Missing data always becomes the main obstacles for the researchers to further their studies. Some researcher will just ignore, truncate, censor, or collapse with those missing data. This might able to make the problem easier but it will lead to inappropriate conclusion and confusion. Therefore, a proper strategy should be used to treat such missing data.

1.2 OBJECTIVE OF THE STUDY

This research is carried out with some objectives as listed below:

- 1) To apply the EM algorithm on multinomial model in missing categorical data analysis.
- 2) To compare the results of independence test for complete and incomplete data.

1.3 SCOPE OF THE STUDY

This study is concentrated on the contingency table where some missing values are present and thus the EM algorithm will be applied on it. Only Missing At Random (MAR) data and Not Missing At Random (NMAR) data are considered in this study.

1.4 SIGNIFICANCE OF THE STUDY

The EM algorithm will be successful in dealing with missing data values in contingency table or in other words we can say that we can find the missing values by applying the EM algorithm. By the end of this study, we will discover a new dimension of problem such as the missingness mechanism which will have a direct impact or effect on the missing values.